

A Coding Theorem for a Class of Stationary Channels with Feedback

Young-Han Kim
University of California, San Diego

February 1, 2008

Abstract

A coding theorem is proved for a class of stationary channels with feedback in which the output $Y_n = f(X_{n-m}^n, Z_{n-m}^n)$ is the function of the current and past m symbols from the channel input X_n and the stationary ergodic channel noise Z_n . In particular, it is shown that the feedback capacity is equal to

$$\lim_{n \rightarrow \infty} \sup_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n),$$

where $I(X^n \rightarrow Y^n) = \sum_{i=1}^n I(X^i; Y_i | Y^{i-1})$ denotes the Massey directed information from the channel input to the output, and the supremum is taken over all causally conditioned distributions $p(x^n || y^{n-1}) = \prod_{i=1}^n p(x_i | x^{i-1}, y^{i-1})$. The main ideas of the proof are the Shannon strategy for coding with side information and a new elementary coding technique for the given channel model *without feedback*, which is in a sense dual to Gallager's lossy coding of stationary ergodic sources. A similar approach gives a simple alternative proof of coding theorems for finite state channels by Yang–Kavčić–Tatikonda, Chen–Berger, and Permuter–Weissman–Goldsmith.

1 Introduction

Shannon [34] showed that the capacity C of a memoryless channel $(\mathcal{X}, p(y|x), \mathcal{Y})$, operationally defined as supremum of all achievable rates [9, Section 7.5], is characterized by

$$C = \sup_{p(x)} I(X; Y). \quad (1)$$

When the channel has memory but still maintains certain ergodic properties, then (1) can be extended to the following multi-letter expression:

$$C = \lim_{n \rightarrow \infty} \sup_{p(x^n)} \frac{1}{n} I(X^n; Y^n). \quad (2)$$

For example, Dobrushin [10] showed that the capacity formula (2) holds if the channel is *information stable*; see also Pinsker [33]. Further extensions and refinements of (2) with more general capacity formulas abound in the literature. For stationary channels, readers are referred to Gray and Ornstein [17], Kieffer [20], and the references therein. A general formula for the capacity is given

by Verdú and Han [38] for arbitrary nonstationary channels that can be represented through a sequence of n -dimensional conditional distributions (even without any consistency requirement); see also Han [18].

For memoryless channels *with feedback*, it was again Shannon [35] who showed that feedback does not increase the capacity and hence that the feedback capacity is given by

$$C_{FB} = C = \sup_{p(x)} I(X; Y). \quad (3)$$

As in the case of nonfeedback capacity (2), the question arises how to extend the feedback capacity formula (3) to channels with memory. The most natural candidate is the following multi-letter expression with *directed information* introduced by Massey [26] in place of the usual mutual information in (2):

$$\begin{aligned} C_{FB} &= \lim_{n \rightarrow \infty} \sup_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n) \\ &= \lim_{n \rightarrow \infty} \sup_{p(x^n || y^{n-1})} \frac{1}{n} \sum_{i=1}^n I(X^i; Y_i | Y^{i-1}), \end{aligned} \quad (4)$$

where the supremum is taken over all n -dimensional *causally conditioned* probabilities

$$\begin{aligned} p(x^n || y^{n-1}) &= \prod_{i=1}^n p(x_i | x^{i-1}, y^{i-1}) \\ &= p(x_1) p(x_2 | x_1, y_1) \cdots p(x_n | x_1, \dots, x_{n-1}, y_1, \dots, y_{n-1}). \end{aligned}$$

The main goal of this paper is to establish the validity of the feedback capacity formula (4) for a reasonably general class of channels with memory, in the simplest manner.

Massey [26] introduced the mathematical notion of directed information

$$I(X^n \rightarrow Y^n) = \sum_{i=1}^n I(X^i; Y_i | Y^{i-1}),$$

and established its operational meaning by showing that the feedback capacity is upper bounded by the maximum normalized directed information, which can be in general tighter than the usual mutual information. He also showed that (4) reduces to (3) if the channel is memoryless, and to (2) if the channel is used without feedback. Kramer [23, 24] streamlined the notion of directed information further and explored many interesting properties; see also Massey and Massey [27].

For channels with certain structures, the validity of the feedback capacity formula (4) has been established implicitly. For example, Cover and Pombra [8] gives a multi-letter characterization of the Gaussian feedback capacity, and Alajaji [1] characterizes the feedback capacity of discrete channels with additive noise—feedback does not increase the capacity of discrete additive channels when there is no input cost constraint. Both results can be recast in the form of directed information (see [8, Eq. (52)] and [1, Eq. (17)]). The notion of directed information in these contexts, however, has a very limited role as an intermediate step in the proof of converse coding theorems. Indeed, the highlight of Cover–Pombra characterization is the asymptotic equipartition property of arbitrary

nonstationary nonergodic Gaussian processes [8, Section V]; see also Pinsker [33]. (The case of discrete additive channel is trivial since the optimal input distribution is memoryless and uniform.)

In a heroic effort [37], Tatikonda attacked the general nonanticipatory channel with feedback by combining Verdú–Han formula for nonfeedback capacity, Massey directed information, and Shannon strategy for channel side information [36]. As the cost of generality, however, it is extremely difficult to establish a simple formula like (4). Furthermore, the coding theorem in [37] is not proved in a completely satisfactory manner.

More recently, Yang, Kavčić, and Tatikonda [40] and Chen and Berger [6] studied special cases of finite-state channels, based on Tatikonda’s framework. A finite-state channel [14, Section 4.6] is described by a conditional probability distribution

$$p(y_n, s_n | x_n s_{n-1}), \quad (5)$$

where s_n denotes the channel state at time n . Using a different approach based on Gallager’s proof of the nonfeedback capacity [14, Section 5.9], Permuter, Weissman, and Goldsmith [31] proved various coding theorems for finite-state channels with feedback that include *inter alia* the results of [40, 6] and establish the validity of (4) for indecomposable finite-state channels without intersymbol interference (i.e., the channel states evolve as an ergodic Markov chain, independent of the channel input).

As mentioned before, we strive to give a straightforward treatment of the feedback coding theorem. Towards this goal, this paper focuses on stationary nonanticipatory channels of the form

$$Y_n = g(X_{n-m}, X_{n-m+1}, \dots, X_n, Z_{n-m}, Z_{n-m+1}, \dots, Z_n). \quad (6)$$

In words, the channel output Y_n at time n is given as a time-invariant deterministic function of channel inputs $X_{n-m}^n = (X_{n-m}, X_{n-m+1}, \dots, X_n)$ up to past m symbols and channel noises $Z_{n-m}^n = (Z_{n-m}, Z_{n-m+1}, \dots, Z_n)$ up to past m symbols. We assume the noise process $\{Z_n\}_{n=1}^\infty$ is an arbitrary stationary ergodic process (without any mixing condition) independent of the message sent over the channel.

The channel model (6) is rather simple and physically motivated. Yet this channel model is general enough to include many important feedback communication models such as any additive noise fading channels with intersymbol interference and indecomposable finite-state channels without intersymbol interference.¹

The channel (6) has finite input memory in the sense of Feinstein [11] and can be viewed as a finite-window sliding-block coder [16, Section 9.4] of input and noise processes (cf. primitive channels introduced by Neuhoff and Shields [29] in which the noise process is memoryless). Compared to the general finite-state channel model (5) in which the channel has infinite input memory but the channel noise is memoryless, our channel model (6) has finite input memory but the noise has infinite memory; recall that there is no mixing condition on the noise process $\{Z_n\}_{n=1}^\infty$. Thus, the finite-state channel model and the finite sliding-block channel model nicely complement each other.

Our main result is to show that the feedback capacity C_{FB} of the channel (6) is characterized by (4). More precisely, we consider a communication problem depicted in Figure 1. Here one wishes

¹A notable exception is a famous finite-state channel called the “trapdoor channel” introduced by Blackwell [3], the feedback capacity of which is established in [30].

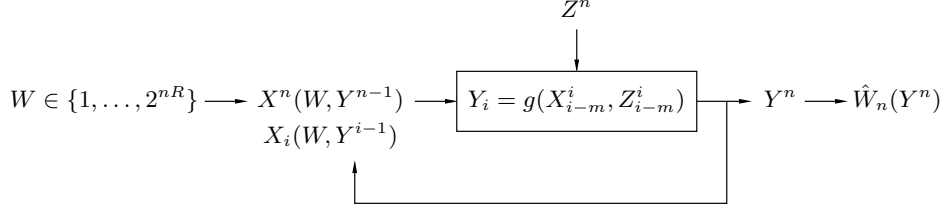


Figure 1: Feedback communication channel $Y_i = g(X_{i-m}^i, Z_{i-m}^i)$.

to communicated a message index $W \in \{1, 2, \dots, 2^{nR}\}$ over the channel

$$Y_i = \begin{cases} \emptyset, & i = 1, \dots, m, \\ g(X_{i-m}^i, Z_{i-m}^i), & i = m+1, m+2, \dots, \end{cases} \quad (7)$$

where the time- i channel output Y_i on the output alphabet \mathcal{Y} is given by a deterministic map $f : \mathcal{X}^m \times \mathcal{Z}^m \rightarrow \mathcal{Y}$ of the current and past m channel inputs X_{i-m}^i on the input alphabet \mathcal{X} and the current and past m channel noises Z_{i-m}^i on the noise alphabet \mathcal{Z} . We assume that the channel noise process $\{Z_i\}_{i=1}^\infty$ is stationary ergodic and is independent of the message W . The initial values of Y_1, \dots, Y_m are set arbitrarily. They depend on the unspecified initial condition (X_{-m+1}^0, Z_{-m+1}^0) , the effect of which vanishes from time $m+1$. Thus the long term behavior of the channel is independent of Y_1^m .

We specify a $(2^{nR}, n)$ feedback code with the encoding maps

$$X^n(W, Y^{n-1}) = (X_1(W), X_2(W, Y_1), \dots, X_n(W, Y^{n-1})), \quad W = 1, \dots, 2^{nR},$$

and the decoding map

$$\hat{W}_n : \mathcal{Y}^n \rightarrow \{1, \dots, 2^{nR}\}.$$

The probability of error $P_e^{(n)}$ is defined as

$$\begin{aligned} P_e^{(n)} &= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \Pr\{\hat{W}_n(Y^n) \neq w | X^n = X^n(w, Y^{n-1})\} \\ &= \Pr\{\hat{W}_n(Y^n) \neq W\}, \end{aligned}$$

where the message W is uniformly distributed over $\{1, \dots, 2^{nR}\}$ and is independent of $\{Z_i\}_{i=1}^\infty$. We say that the rate R is achievable if there exists a sequence of $(2^{nR}, n)$ codes with $P_e^{(n)} \rightarrow 0$ as $n \rightarrow \infty$. The feedback capacity C_{FB} is defined as the supremum of all achievable rates. The nonfeedback capacity C is defined similarly, with codewords $X^n(W) = (X_1(W), \dots, X_n(W))$ restricted to be a function of the message W only.

We will prove the following result in Section 4.

Theorem 1. *The feedback capacity C_{FB} of the channel (7) is given by*

$$C_{FB} = \lim_{n \rightarrow \infty} \sup_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n). \quad (8)$$

Our development has two major ingredients. First, we revisit the communication problem over the same channel without feedback in Section 3 and prove that the nonfeedback capacity is given by

$$C = \lim_{n \rightarrow \infty} \sup_{p(x^n)} \frac{1}{n} I(X^n; Y^n).$$

Roughly speaking, there are three flavors in the literature for the achievability proof of non-feedback capacity theorems. The first one is Shannon’s original argument [34] based on random codebook generation, asymptotic equipartition property, and joint typicality decoding, which was made rigorous by Forney [13] and Cover [7], and now is used widely in coding theorems for memoryless networks [9, Chapter 15]. This approach, however, does not easily generalize to channels with memory. The second flavor is the method of random coding exponent by Gallager [15], which was later applied to finite-state channels [14, Section 5.9]. This approach is perhaps the simplest one for the analysis of general finite-state channels and has been adapted by Lapidoth and Telatar [25] for compound finite-state channels and by Permuter *et al.* [31] for finite-state channels with feedback.

The third and the least intuitive approach is Feinstein’s fundamental lemma [12]. This is the most powerful and general method of the three, and has been applied extensively in the literature, say, from Khinchin [19] to Gray [16] to Verdú and Han [38].

Our approach is somewhat different from these three usual approaches. We use the strong typicality (relative frequency) decoding for n -dimensional super letters. A constructive coding scheme (up to the level of Shannon’s random codebook generation) based on block ergodic decomposition of Nedoma [28] is developed, which uses a long codeword on the n -letter super alphabet, constructed as a concatenation of n shorter codewords. While each short codeword and the corresponding output fall into their own ergodic mode, the long codeword as a whole maintains the ergodic behavior. To be fair, codebook construction of this type is far from new in the literature, and our method is intimately related to the one used by Gallager [14, Section 9.8] and Berger [2, Section 7.2] for lossy compression of stationary ergodic sources. Indeed, when the channel (6) has zero memory ($m = 0$), then the role of the input for our channel coding scheme is equivalent to the role of the covering channel for Gallager’s source coding scheme.

Equipped with this coding method for nonfeedback sliding-block coder channels (6), the extension to the feedback case is relatively straightforward. The basic ingredient for this extension is the Shannon strategy for channels with causal side information at the transmitter [36]. As a matter of fact, Shannon himself observed that the major utility of his result is feedback communication. Following is the first sentence of [36]:

Channels with feedback from the receiving to the transmitting point are a special case of a situation in which there is additional information available at the transmitter which may be used as an aid in the forward transmission system.

As observed by Caire and Shamai [5, Proposition 1], the causality has no cost when the transmitter and the receiver share the same side information—in our case, the past input (if decoded faithfully) and the past output (received from feedback)—and the transmission can fully utilize this side information as if it were known *a priori*.

Intuitively speaking, we can achieve the rate R_i for the i th symbol in the length- n super symbol

as

$$R_i = \max_{p(x^i||y^{i-1})} I(X_i; Y_i^n | X^{i-1}, Y^{i-1}), \quad i = 1, \dots, n,$$

and hence the total achievable rate becomes

$$R = \max_{p(x^n||y^{n-1})} \sum_{i=1}^n I(X_i; Y_i^n | X^{i-1}, Y^{i-1})$$

per n transmissions. Now a simple algebra shows that this rate is equal to the maximum directed information as follows:

$$\begin{aligned} \sum_{i=1}^n I(X_i; Y_i^n | X^{i-1}, Y^{i-1}) &= \sum_{i=1}^n \sum_{j=i}^n I(X_i; Y_j | X^{i-1}, Y^{j-1}) \\ &= \sum_{j=1}^n \sum_{i=1}^j I(X_i; Y_j | X^{i-1}, Y^{j-1}) \\ &= \sum_{j=1}^n I(X^j; Y_j | Y^{j-1}) \\ &= I(X^n \rightarrow Y^n). \end{aligned} \tag{9}$$

The above argument, while intuitively appealing, is not completely rigorous, however. Therefore, we will take more careful steps, by first proving the achievability of $\frac{1}{n}I(U^n; Y^n)$ for all auxiliary random variables U^n and Shannon strategies $X_i(U_i, X^{i-1}, Y^{i-1})$, $i = 1, \dots, n$, and then showing that $I(U^n; Y^n)$ reduces to $I(X^n \rightarrow Y^n)$ via pure algebra.

The next section collects all necessary lemmas that will be used subsequently in Section 3 for the nonfeedback coding theorem and in Section 4 for the feedback coding theorem.

2 Preliminaries

Here we review relevant materials from ergodic theory and information theory in the form of 10 lemmas. While some of the lemmas are classical and are presented in order to make the paper self-contained, the other lemmas are crucial to our main discussion in subsequent sections and may contain original observations. Throughout this section, $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$ denotes a generic stochastic process on a finite alphabet \mathcal{Z} with associated probability measure P defined on Borel sets under the usual topology on \mathcal{Z}^∞ .

2.1 Ergodicity

Given a *stationary* process $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$, let $T : \mathcal{Z}^\infty \rightarrow \mathcal{Z}^\infty$ be the associated measure preserving shift transformation. Intuitively, T maps the infinite sequence (z_1, z_2, z_3, \dots) to (z_2, z_3, z_4, \dots) . We say the transformation T (or the process \mathbf{Z} itself) is *ergodic* if every measurable set A with $TA = A$ satisfies either $P(A) = 0$ or $P(A) = 1$.

The following characterization of ergodicity is well known; see, for example, Petersen [32, Exercise 2.4.4] or Wolfowitz [39, Lemma 10.3.1].

Lemma 1. Suppose $\{Z_i\}_{i=1}^\infty$ be a stationary process and let T denote the associated measure preserving shift transformation. Then, $\{Z_i\}$ is ergodic if and only if

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} P(T^{-i}A \cap B) = P(A) \cdot P(B) \quad \text{for all measurable } A \text{ and } B.$$

When $\mathbf{X} = \{X_i\}_{i=1}^\infty$ and $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$ are independent stationary ergodic processes, they are not necessarily jointly ergodic. For example, if we take

$$\mathbf{X} = \begin{cases} 01010101 \dots, & \text{with probability } 1/2, \\ 10101010 \dots, & \text{with probability } 1/2, \end{cases}$$

and \mathbf{Z} is independent and identically distributed as \mathbf{X} , then it is easy to verify that $\{Y_i = X_i + Z_i \pmod{2}\}_{i=1}^\infty$ is *not* ergodic. However, if one of the processes is mixing reasonably fast, then they are jointly ergodic. The following result states a sufficient condition for joint ergodicity.

Lemma 2. If \mathbf{X} is independent and identically distributed (i.i.d.), and \mathbf{Z} is stationary ergodic, independent of \mathbf{X} , then the pair $(\mathbf{X}, \mathbf{Z}) = \{(X_i, Z_i)\}_{i=1}^\infty$ is jointly stationary ergodic.

A stronger result is true, which assumes \mathbf{X} to be weakly mixing only. The proof is an easy consequence of Lemma 1; for details refer to Brown [4, Proposition 1.6] or Wolfowitz [39, Theorem 10.3.1].

We will later need to construct *super-letter* processes for our coding theorems. The next lemma due to Gallager [14, Lemma 9.8.2] deals with the ergodic decomposition of the n -letter super process that is built from a single-letter stationary ergodic one; see also Nedoma [28] and Berger [2, Section 7.2].

Lemma 3. Suppose $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$ be stationary ergodic on \mathcal{Z} , and let T be the associated shift transformation. Define the n th-order super process $\mathbf{Z}^{(n)} = \{Z_i^{(n)}\}_{i=1}^\infty$ on \mathcal{Z}^n as

$$Z_i^{(n)} = (Z_{n(i-1)+1}, Z_{n(i-1)+2}, \dots, Z_{ni}), \quad i = 1, 2, \dots$$

Then, the super process $\mathbf{Z}^{(n)}$ has n' ergodic modes, each with probability $1/n'$ and disjoint up to measure zero, where n' divides n . Furthermore, in the space \mathcal{Z}^∞ of the original process \mathbf{Z} , the sets $S_1, S_2, \dots, S_{n'}$ corresponding to these ergodic modes can be related by $T(S_i) = S_{i+1}$, $1 \leq i \leq n-1$, and $T(S_n) = S_1$.

We will use the notation $P(\cdot|S_k)$, $k = 1, \dots, n'$, for the probability measure under each ergodic mode.

2.2 Strong Typicality

We use the strong typicality [9, Section 10.6] as the basic method of decoding. Here we review a few basic properties of strongly typical sequences.

First definitions. Let $N(a|x^n)$ denote the number of occurrences of the symbol a in the sequence x^n . We say a sequence $x^n \in \mathcal{X}^n$ is ϵ -strongly typical (or typical in short) with respect to a

distribution $P(x)$ on \mathcal{X} if

$$\left| \frac{1}{n} N(a|x^n) - P(a) \right| < \frac{\epsilon}{|\mathcal{X}|}$$

for all $a \in \mathcal{X}$ with $P(a) > 0$, and $N(a|x^n) = 0$ for all $a \in \mathcal{X}$ with $P(a) = 0$. Consistent with this definition, we say a pair of sequences (x^n, y^n) are *jointly ϵ -strongly typical* (or jointly typical in short) with respect to a distribution $P(x, y)$ on $\mathcal{X} \times \mathcal{Y}$ if

$$\left| \frac{1}{n} N(a, b|x^n, y^n) - P(a, b) \right| < \frac{\epsilon}{|\mathcal{X}||\mathcal{Y}|}$$

for all $(a, b) \in \mathcal{X} \times \mathcal{Y}$ with $P(a, b) > 0$, and $N(a, b|x^n, y^n) = 0$ for all $(a, b) \in \mathcal{X} \times \mathcal{Y}$ with $P(a, b) = 0$.

The set of strongly typical sequences $x^n \in \mathcal{X}^n$ with respect to $X \sim P(x)$ is denoted $A_\epsilon^{*(n)}(X)$. We similarly define a joint typical set $A_\epsilon^{*(n)}(X, Y)$ for $(X, Y) \sim P(x, y)$.

The following statement is a trivial consequence of the definition of typical sequences.

Lemma 4. *Suppose $X \sim P(X)$. If $x^n \in A_\epsilon^{*(n)}(X)$ and $y^n = f(x^n) := (f(x_1), f(x_2), \dots, f(x_n))$, then $y^n \in A_\delta^{*(n)}(f(X))$ with $\delta = \epsilon \cdot (|\mathcal{X}| - 1)$.*

As a special case, if (x^n, y^n) is ϵ -strongly typical with respect to a joint distribution $P(x, y)$, then x^n is ϵ -strongly typical with respect to the marginal $P(x) = \sum_y P(x, y)$.

Our discussion on the typical sequences so far has not given a specific context on how they are generated. Now we connect the notion of strong typicality with ergodic processes. First, from Birkhoff's ergodic theorem [32, Theorem 2.2.3] and the definition of ergodicity, the following lemma is immediate.

Lemma 5. *Let $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$ be stationary ergodic with $Z_1 \sim P(z)$. Then*

$$\Pr(Z^n \in A_\epsilon^{*(n)}(Z_1)) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

As we mentioned in the previous subsection, the n th order super process $\mathbf{Z}^{(n)} = \{Z_i^{(n)}\}_{i=1}^\infty$ defined as

$$Z_i^{(n)} = Z_{(n-1)i+1}^{ni}, \quad i = 1, 2, \dots,$$

is not necessarily ergodic, but is a mixture of disjoint ergodic modes. Thus, the super process $\mathbf{Z}^{(n)}$ is not necessarily typical with respect to $P(z^n)$ on the n -letter alphabet \mathcal{Z}^n . The following construction by Gallager [14, pp. 498–499], however, gives a typical sequence in the n -letter super alphabet by shifting through each ergodic phase.

Lemma 6. *Given positive integers n, L and a stationary ergodic process $\mathbf{Z} = \{Z_i\}_{i=1}^\infty$, construct $\tilde{\mathbf{Z}} = \{\tilde{Z}_i\}_{i=1}^{Ln^2}$ as follows (See Figure 2):*

$$\tilde{Z}_i = \begin{cases} Z_i, & i = 1, \dots, Ln, \\ Z_{i+1}, & i = Ln + 1, \dots, 2Ln, \\ \vdots & \\ Z_{i+n-1}, & i = Ln(n-1) + 1, \dots, Ln^2. \end{cases}$$

In other words, $\{\tilde{Z}_i\}_{i=1}^{Ln^2}$ is a verbatim copy of $\{Z_i\}_{i=1}^{Ln^2+n}$ with every $(Ln+1)$ st position skipped.

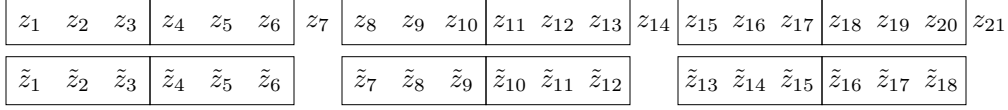


Figure 2: Construction of \tilde{Z}^{Ln^2} from Z^{Ln^2+n} : $n = 3, L = 2$.

Let $\tilde{\mathbf{Z}}^{(n)} = (\tilde{Z}_1^n, \tilde{Z}_{n+1}^{2n}, \dots, \tilde{Z}_{Ln^2-n+1}^{Ln^2})$ be the associated n th order super process of length Ln . Then,

$$\Pr(\tilde{\mathbf{Z}}^{(n)} \in A_\epsilon^{*(Ln)}(Z^n)) \rightarrow 1 \quad \text{as } L \rightarrow \infty.$$

Proof. From Lemma 3 and the given construction of skipping one position after every Ln symbols, each of n sequences

$$\begin{aligned}
(\tilde{Z}_1^n, \dots, \tilde{Z}_{Ln-n+1}^{Ln}) &= (Z_1^n, \dots, Z_{Ln-n+1}^{Ln}) \\
(\tilde{Z}_{Ln+1}^{Ln+n}, \dots, \tilde{Z}_{2Ln-n+1}^{2Ln}) &= (Z_{Ln+2}^{L(n+1)+1}, \dots, Z_{2Ln-n+2}^{2Ln+1}) \\
&\vdots \\
(\tilde{Z}_{Ln(n-1)+1}^{Ln(n-1)+n}, \dots, \tilde{Z}_{Ln^2-n+1}^{Ln^2}) &= (Z_{Ln(n-1)+n}^{Ln(n-1)+2n-1}, \dots, Z_{Ln^2-n+1}^{Ln^2+n-1})
\end{aligned}$$

falls in one of ergodic modes $(S_1, \dots, S_{n'})$ with n/n' sequences for each mode. Now for each sequence with corresponding ergodic mode S_k , the relative frequencies of all super symbols $a^n \in \mathcal{Z}^n$ converge to the corresponding distribution $P(a^n|S_k)$ as $L \rightarrow \infty$. But each ergodic mode is visited evenly, each by n/n' sequences. Therefore, the relative frequencies of all $a^n \in \mathcal{Z}^n$ in the entire sequence $\tilde{Z}_1^{Ln^2}$ converge to

$$\frac{1}{n'} \sum_{k=1}^{n'} P(a^n|S_k) = P(a^n)$$

as $L \rightarrow \infty$. □

Combining Lemma 2 with the proof of Lemma 6, we have the following result.

Lemma 7. *Under the condition of Lemma 6, let further $\mathbf{X} = \{X_i\}_{i=1}^\infty$ be blockwise i.i.d. $\sim P(x^n)$, that is, $X_i^{(n)} = X_{(n-1)i+1}^{ni}$, $i = 1, 2, \dots$, i.i.d. $\sim P(x^n)$, independent of \mathbf{Z} . Then,*

$$\Pr((\mathbf{X}^{(n)}, \tilde{\mathbf{Z}}^{(n)}) \in A_\epsilon^{*(Ln)}(X^n, Z^n)) \rightarrow 1 \quad \text{as } L \rightarrow \infty.$$

Finally we recall the key result linking the typicality with mutual information [9, Lemma 10.6.2].

Lemma 8. *Suppose $(X, Y) \sim P(x, y)$ and let X_1, X_2, \dots, X_n be i.i.d. $\sim P(x)$. For $y^n \in A_\epsilon^{*(n)}(Y)$, the probability that $(X^n, y^n) \in A_\epsilon^{*(n)}(X, Y)$ is upper bounded by*

$$\Pr((X^n, y^n) \in A_\epsilon^{*(n)}(X, Y)) \leq 2^{-n(I(X;Y)-\delta)}$$

where $\delta \rightarrow 0$ as $\epsilon \rightarrow 0$.

2.3 Channels with Side Information

We prove the following identity in a purely algebraic manner, then find its meaning in information theory. Here we assume every alphabet is finite.

Lemma 9. *Suppose $S \sim p(s)$. For a given conditional distribution $p(y|x, s)$ on the product space $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$, we have*

$$\max_{p(u), x=f(u,s)} I(U; Y, S) = \max_{p(x|s)} I(X; Y|S), \quad (10)$$

where the maximum on the left hand side is taken over all conditional distributions of the form $p(u, x|s) = p(u)p(x|u, s)$ with deterministic $p(x|u, s)$ (that is, $p(x|u, s) = 0$ or 1), and the auxiliary random variable U has cardinality bounded by $|\mathcal{U}| \leq (|\mathcal{X}| - 1)|\mathcal{S}|$.

Proof. For any joint distribution of the form $p(u, x, s, y) = p(u)p(s)p(x|u, s)p(y|x, s)$ with deterministic $p(x|u, s)$, we have the following Markov chains: $U \rightarrow (X, S) \rightarrow Y$ and $X \rightarrow (U, S) \rightarrow Y$. Combined with the independence of U and S , these Markov relationships imply that

$$\begin{aligned} \max_{p(u), x=f(u,s)} I(U; Y, S) &= \max_{p(u), x=f(u,s)} I(U; Y|S) \\ &= \max_{p(u), x=f(u,s)} I(X; Y|S). \end{aligned} \quad (11)$$

But it can be easily verified that any conditional distribution $p(x|s)$ can be represented as

$$p(x|s) = \sum_u p(u)p(x|u, s)$$

for appropriately chosen $p(u)$ and *deterministic* $p(x|u, s)$ with cardinality of U upper bounded by $|\mathcal{U}| \leq (|\mathcal{X}| - 1)|\mathcal{S}|$. Therefore, we have

$$\max_{p(u), x=f(u,s)} I(X; Y|S) = \max_{p(x|s)} I(X; Y|S),$$

which proves the desired result. \square

It is well known that the capacity of a memoryless state-dependent channel $p(y|x, s)$ is given as

$$C = \max_{p(x|s)} I(X; Y|S),$$

if the state information is known at both the encoder and decoder prior to the actual communication. What will happen if the transmitter learns the state information on the fly, so that only the past and present state realization can be utilized for communication?

Shannon [36] considered the communication over a memoryless state-dependent channel $p(y|x, s)$ with state information available *only* at the transmitter on the fly, and showed that the capacity is given by

$$C = \max_{p(u), x=f(u,s)} I(U; Y), \quad (12)$$

where the cardinality of U is bounded as $|\mathcal{U}| \leq |\mathcal{X}|^{|\mathcal{S}|}$, counting for all functions $f : \mathcal{S} \rightarrow \mathcal{X}$. This capacity is achieved by attaching a physical device $X = f(U, S)$ in front of the actual channel as

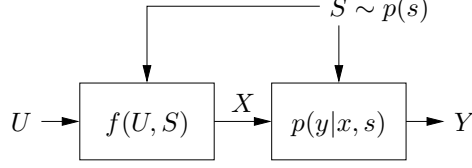


Figure 3: Shannon strategy for coding with side information.

depicted in Figure 3, which maps the channel state S to the channel input X according to the function (index) U . Now treating U as the input to the newly generated channel

$$p(y|u) = \sum_{x,s} p(s)p(x|u,s)p(y|x,s)$$

and coding as in the case of usual memoryless channels, we can easily achieve $I(U; Y)$. This method, surprisingly simple yet optimal, is sometimes called the Shannon strategy.

Now when the decoder also knows the channel state S , it is equivalent for the decoder to receive the augmented channel output $Y' = (Y, S)$. Thus, the capacity of the same channel $p(y|x, s)$ with the state information *causally* known at both the encoder and decoder² follows from (12) as

$$C = \max_{p(u), x=f(u,s)} I(U; Y, S).$$

Therefore, Lemma 9 states that when the same side information is available at the receiver, the causal encoder with the best Shannon strategy performs no worse than the noncausal encoder who can preselect the entire codeword compatible with the whole state sequence.

For the last lemma needed for main results, we recall the notation of causally conditioned distributions

$$p(x^n || y^{n-1}) = \prod_{i=1}^n p(x_i | x^{i-1}, y^{i-1}) \quad (13)$$

and

$$p(y^n || x^n) = \prod_{i=1}^n p(y_i | x^i, y^{i-1}). \quad (14)$$

(The notation (13) and (14) can be unified if we define

$$p(a^n || b^m) = \begin{cases} \prod_{i=1}^n p(a_i | b^i, a^{i-1}), & n = m, \\ p(a^n || \emptyset^{n-m} b^m), & n > m, \\ p(\emptyset^{m-n} a^n || b^m), & n < m. \end{cases}$$

²For the usual block coding, the decoder causality is irrelevant. The message is decoded only after the entire block is received.

By chain rule, we have

$$p(x^n||y^{n-1})p(y^n||x^n) = p(x^n, y^n) = p(x^n)p(y^n|x^n)$$

for any joint distribution $p(x^n, y^n)$. Thus, given a causally conditioned distribution (or a channel) $p(y^n||x^n)$, the causally conditioned distribution (or the input) $p(x^n||y^{n-1})$ completely specifies the joint distribution $p(x^n, y^n)$.

As a corollary of Lemma 9, we have the following result.

Lemma 10. *Suppose a causally conditioned distribution $p(y^n||x^n)$ is given. Then we have*

$$\max_{p(u^n), x_i=f(u_i, x^{i-1}, y^{i-1})} I(U^n; Y^n) = \max_{p(x^n||y^{n-1})} I(X^n \rightarrow Y^n), \quad (15)$$

where the maximum on the left hand side is taken over all joint distributions of the form

$$\begin{aligned} p(u^n, x^n, y^n) &= \prod_{i=1}^n (p(u_i)p(x_i|u_i, x^{i-1}, y^{i-1})p(y_i|x^i, y^{i-1})) \\ &= \left(\prod_{i=1}^n p(u_i)p(x_i|u_i, x^{i-1}, y^{i-1}) \right) p(y^n||x^n) \end{aligned} \quad (16)$$

with deterministic $p(x_i|u_i, x^{i-1}, y^{i-1})$, $i = 1, \dots, n$, and the auxiliary random variables U_i has the cardinality bounded by $|\mathcal{U}_i| \leq |\mathcal{X}|^i |\mathcal{Y}|^{i-1}$.

Proof. Let $q(u^n, x^n, y^n)$ be any joint distribution of the form (16) such that $q(x_i|u_i, x^{i-1}, y^{i-1})$, $i = 1, \dots, n$ are deterministic and that $q(y^n||x^n) = p(y^n||x^n)$ (i.e., the joint distribution $q(u^n, x^n, y^n)$ is consistent with the given causally conditioned distribution $p(y^n||x^n)$). For $(U^n, X^n, Y^n) \sim q(u^n, x^n, y^n)$, it is easy to verify that U_i^n is independent of $(U^{i-1}, X^{i-1}, Y^{i-1})$, which implies that $U^{i-1} \rightarrow (X^{i-1}, Y^{i-1}) \rightarrow Y_i^n$ forms a Markov chain. On the other hand, X^{i-1} is a deterministic function of (U^{i-1}, Y^{i-1}) and thus $X^{i-1} \rightarrow (U^{i-1}, Y^{i-1}) \rightarrow Y_i^n$ also forms a Markov chain. Similarly, we have the Markovity for $U^i \rightarrow (X^i, Y^{i-1}) \rightarrow Y_i^n$ and $X^i \rightarrow (U^i, Y^{i-1}) \rightarrow Y_i^n$. Therefore, we have

$$I(U_i; Y^n | U^{i-1}) = I(U_i; Y_i^n | Y^{i-1}, U^{i-1}) \quad (17)$$

$$\begin{aligned} &= H(Y_i^n | Y^{i-1}, U^{i-1}) - H(Y_i^n | Y^{i-1}, U^i) \\ &= H(Y_i^n | Y^{i-1}, X^{i-1}) - H(Y_i^n | Y^{i-1}, X^i) \\ &= I(X_i; Y_i^n | X^{i-1}, Y^{i-1}), \end{aligned} \quad (18)$$

where (17) follows from the independence of U_i and (U^{i-1}, Y^{i-1}) , and (18) follows from Markov relationships observed above. Now from the alternative expansion of the directed information shown in (9), we have

$$\max_q I(U^n; Y^n) = \max_q I(X^n \rightarrow Y^n).$$

Finally, by using distributions of the form

$$p(x_i|x^{i-1}, y^{i-1}) = \sum_{u_i} p(u_i)p(x_i|u_i, x^{i-1}, y^{i-1}), \quad i = 1, \dots, n$$

with appropriately chosen $p(u_i)$ and deterministic $p(x_i|u_i, x^{i-1}, y^{i-1})$, we can represent any causally conditioned distribution

$$p(x^n||y^{n-1}) = \prod_{i=1}^n p(x_i|x^{i-1}, y^{i-1}) = \sum_{u^n} \prod_{i=1}^n (p(u_i)p(x_i|u_i, x^{i-1}, y^{i-1})),$$

which implies that

$$\max_q I(X^n \rightarrow Y^n) = \max_{p(x^n||y^{n-1})} I(X^n \rightarrow Y^n)$$

and completes the proof. \square

3 Nonfeedback Coding Theorem Revisited

This section is devoted to the proof of the following result.

Theorem 2. *The nonfeedback capacity C of the stationary channel*

$$Y_i = \begin{cases} \emptyset, & i = 1, \dots, m, \\ g(X_{i-m}^i, Z_{i-m}^i), & i = m+1, m+2, \dots, \end{cases} \quad (19)$$

with the input X_i and the stationary ergodic noise process $\{Z_i\}_{i=1}^\infty$ depicted in Figure 1 is given by

$$\begin{aligned} C &= \lim_{n \rightarrow \infty} C_n \\ &= \lim_{n \rightarrow \infty} \sup_{p(x^n)} \frac{1}{n} I(X^n; Y^n). \end{aligned} \quad (20)$$

Revisiting and proving the nonfeedback coding theorem is rewarding for two reasons. First, our proof is somewhat different from the usual techniques and hence is interesting on its own. (See Section 1 for the discussion on conventional achievability proofs of nonfeedback capacity theorems.) Second, our exercise here will lead to a straightforward proof of the feedback coding theorem in the next section.

Proof. We first note that the capacity expression (20) is well-defined because nC_n is superadditive (i.e., $mC_m + nC_n \leq (m+n)C_{m+n}$), which implies that the limit exists and

$$\lim_{n \rightarrow \infty} C_n = \sup_{n \geq 1} C_n.$$

The converse follows immediately from Fano's inequality [9, Lemma 7.9.1]. For any sequence of $(2^{nR}, n)$ codes $(X^n(W), \hat{W}(Y^n))$ with the message W drawn uniformly over $\{1, \dots, 2^{nR}\}$, if

$$P_e^{(n)} = \Pr(W \neq \hat{W}) \rightarrow 0,$$

then we must have

$$\begin{aligned} nR &\leq I(W; Y^n) + n\epsilon_n \\ &\leq I(X^n; Y^n) + n\epsilon_n \end{aligned}$$

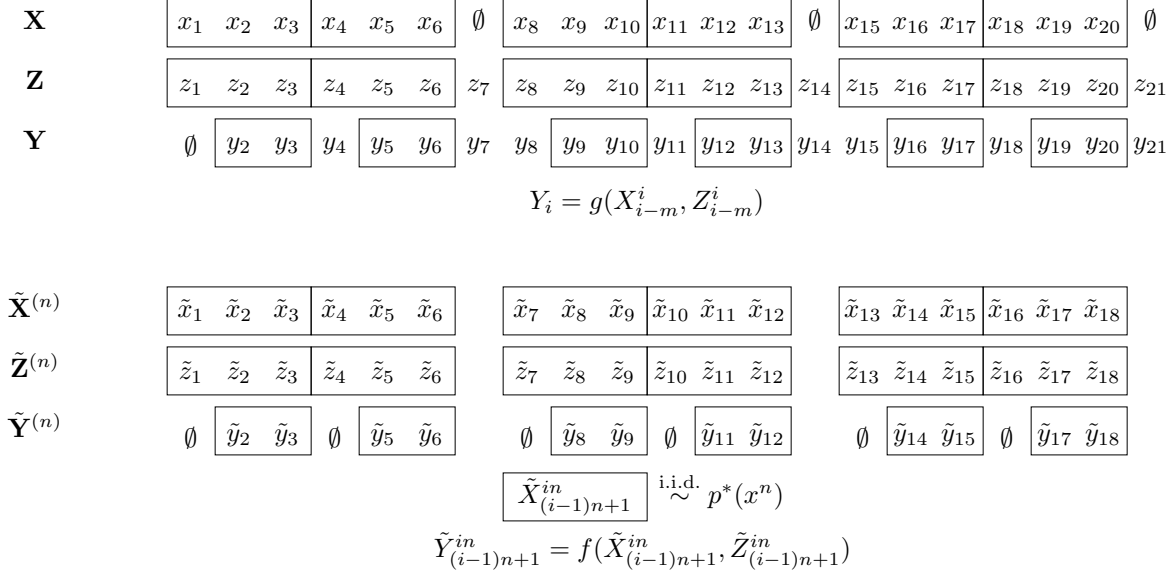


Figure 4: Input, noise, and output sequences: $n = 3, L = 2, m = 1$.

where $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$.

For the achievability, it suffices to show that there exists a sequence of codes that achieves C_n for each $n > m$. (Recall $C_1 = \dots C_m = 0$.) Without loss of generality, we assume that the alphabets $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ are finite. Otherwise, we can partition the space for each n and $\epsilon > 0$ such that

$$\max_{p([x]^n)} \frac{1}{n} I([X]^n; [Y]^n) \geq C_n - \epsilon,$$

and prove the achievability on this partitioned space.

Codebook generation. Fix $n > m$ and let $p^*(x^n)$ denote the input distribution that achieves C_n . For each $L = 1, 2, \dots$, let $k = k(L, n) = Ln^2 + n$. We generate a sequence³ of $(2^{kR}, k)$ codes $X^k(w)$ as depicted in Figure 4.

For each $w \in \{1, 2, \dots, 2^{kR}\}$, generate a codeword $\tilde{\mathbf{X}}^{(n)}(w) = \tilde{X}^{Ln^2}(w)$ of length Ln on the n -letter super alphabet \mathcal{X}^n independently according to

$$p(\tilde{x}^{Ln^2}) = \prod_{i=1}^{Ln} p^*(x_{(n-1)i+1}^{ni}).$$

We exhibit the 2^{kR} codewords as the rows of a matrix:

$$\mathcal{C} = \begin{bmatrix} \tilde{X}_1^n(1) & \tilde{X}_{n+1}^{2n}(1) & \dots & \tilde{X}_{Ln(n-1)+1}^{Ln^2}(1) \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{X}_1^n(2^{kR}) & \tilde{X}_{n+1}^{2n}(2^{kR}) & \dots & \tilde{X}_{Ln(n-1)+1}^{Ln^2}(2^{kR}) \end{bmatrix}.$$

Each entry in this matrix is generated i.i.d. according to $p^*(x^n)$.

³This gives only a *subsequence* of $(2^{kR}, k)$ codes. But we can easily interpolate to $Ln^2 + n < k < (L+1)n^2 + n$ without any rate loss, since $(Ln^2 + n)/((L+1)n^2 + n) \rightarrow 1$ as $L \rightarrow \infty$.

Using the construction as in Lemma 6 (see Figure 4), the *actual* codewords $\mathbf{X}(w) = X^k(w)$, $w = 1, 2, \dots, 2^{nR}$, which will be transmitted over the channel, are generated from $\tilde{\mathbf{X}}^{(n)}(w) = \tilde{X}^{Ln^2}(w)$ as follows:

$$X_{(i-1)Ln+i}^{iLn+i} = (\tilde{X}_{(i-1)Ln+1}^{iLn}, \emptyset), \quad i = 1, 2, \dots, n.$$

In other words, X^k is a verbatim copy of \tilde{X}^{Ln^2} with fixed symbol \emptyset separating the subsequences of length Ln .

Encoding. If $W = w$, the transmitter sends the codeword $\mathbf{X}(w) = X^k(w)$ over the channel.

Decoding. Upon receiving the sequence $\mathbf{Y} = Y^k$, the receiver forms the sequence $\tilde{\mathbf{Y}}^{(n)} = \tilde{Y}^{Ln^2}$ of length Ln in the n -letter super alphabet \mathcal{Y}^n , as depicted in Figure 4:

$$\tilde{Y}_{(i-1)n+1}^{in} = \begin{cases} (\emptyset, Y_{(i-1)n+m+1}^{in}), & i = 1, \dots, L, \\ (\emptyset, Y_{(i-1)n+m+2}^{in+1}), & i = L+1, \dots, 2L, \\ \vdots \\ (\emptyset, Y_{(i-1)n+m+n-1}^{in+n}), & i = L(n-1)+1, \dots, Ln. \end{cases}$$

Now we consider $\tilde{\mathbf{X}}^{(n)} = \tilde{X}^{Ln^2}$ and $\tilde{\mathbf{Y}}^{(n)} = \tilde{Y}^{Ln^2}$ as sequences of length Ln on the super alphabet $\mathcal{X}^n \times \mathcal{Y}^n$. The receiver declares that the message \hat{W} was sent if there is a unique \hat{W} such that

$$(\tilde{\mathbf{X}}^{(n)}(\hat{W}), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}(X^n, Y^n),$$

that is, $(\tilde{\mathbf{X}}^{(n)}(\hat{W}), \tilde{\mathbf{Y}}^{(n)})$ is jointly typical with respect to the joint distribution $p(x^n, y^n)$ specified by $p^*(x^n)p(z^n)$ and the definition of the channel (19). Otherwise, an error is declared.

Analysis of the probability of error. Without loss of generality, we assume $W = 1$ was sent. We define the following events:

$$E_i = \{(\tilde{\mathbf{X}}^{(n)}(1), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}(X^n, Y^n)\}, \quad i \in \{1, 2, \dots, 2^{kR}\},$$

where E_i is the event that the i th codeword and $\tilde{\mathbf{Y}}^{(n)}$ are jointly typical. By Bonferonni's inequality, we have

$$\begin{aligned} \Pr(\hat{W} \neq W) &= \Pr(\hat{W} \neq W | W = 1) \\ &= \Pr(E_1^c \cup E_2 \cup E_3 \cup \dots \cup E_{2^{kR}}) \\ &\leq \Pr(E_1^c) + \sum_{i=2}^{2^{kR}} \Pr(E_i). \end{aligned}$$

In order to bound $\Pr(E_1^c)$, we define $\tilde{\mathbf{Z}}^{(n)}$ as the n th order super process of length Ln on the super alphabet \mathcal{Z}^n constructed from the noise process $\{Z_i\}_{i=1}^\infty$ as in Lemma 6. (See Figure 4.) Since $\tilde{\mathbf{X}}^{(n)}(1)$ is blockwise i.i.d. $\sim p^*(x^n)$ and independent of \mathbf{Z} , we have from Lemma 7

$$\Pr((\tilde{\mathbf{X}}^{(n)}(1), \tilde{\mathbf{Z}}^{(n)}) \in A_\epsilon^{*(Ln)}(X^n, Z^n)) \rightarrow 1 \quad \text{as } L \rightarrow \infty.$$

Furthermore, $\tilde{\mathbf{Y}}^{(n)}$ is the blockwise function of $(\tilde{\mathbf{X}}^{(n)}(1), \tilde{\mathbf{Z}}^{(n)})$, that is,

$$\tilde{Y}_{(i-1)n+1}^{in} = f(\tilde{X}_{(i-1)n+1}^{in}(1), \tilde{Z}_{(i-1)n+1}^{in})$$

with the time-invariant function f induced by the channel function g in (19). Thus by Lemma 4,

$$\Pr((\tilde{\mathbf{X}}^{(n)}(1), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}(X^n, Y^n)) \rightarrow 1 \quad \text{as } L \rightarrow \infty,$$

and

$$\Pr(E_1^c) \leq \epsilon \quad \text{for } L \text{ sufficiently large.}$$

On the other hand, recall that the typicality of $(\tilde{\mathbf{X}}^{(n)}(i), \tilde{\mathbf{Y}}^{(n)})$ implies the typicality of $\tilde{\mathbf{Y}}^{(n)}$ (Lemma 4). Hence, by Lemma 8 we have for each $i \neq 1$

$$\begin{aligned} \Pr(E_i) &= \Pr((\tilde{\mathbf{X}}^{(n)}(i), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}) \\ &= \sum_{\tilde{\mathbf{y}}^{(n)} \in A_\epsilon^{*(Ln)}} \Pr((\tilde{\mathbf{X}}^{(n)}(i), \tilde{\mathbf{y}}^{(n)}) \in A_\epsilon^{*(Ln)}) \leq 2^{-Ln(I(X^n; Y_{m+1}^n) - \delta)}, \end{aligned}$$

where $\delta \rightarrow 0$ as $\epsilon \rightarrow 0$. Consequently,

$$\begin{aligned} \Pr(\hat{W} \neq W) &\leq \Pr(E_1^c) + \sum_{i=2}^{2^{kR}} \Pr(E_i) \\ &\leq \epsilon + 2^{kR} 2^{-Ln(I(X^n; Y_{m+1}^n) - \delta)} \\ &\leq 2\epsilon \end{aligned}$$

if L is sufficiently large and

$$kR < Ln(I(X^n; Y_{m+1}^n) - \delta),$$

or equivalently,

$$R < \frac{Ln^2 + n}{Ln} (I(X^n; Y_{m+1}^n) - \delta).$$

Since ϵ can be made arbitrarily small and $(Ln^2 + n)/(Ln) \rightarrow 1/n$ as $L \rightarrow \infty$, we have a sequence of $(2^{kR}, k)$ codes that achieves

$$R < \frac{1}{n} I(X^n; Y_{m+1}^n) = \frac{1}{n} I(X^n; Y^n) = C_n.$$

□

4 Proof of Theorem 1

Recall our channel model:

$$Y_i = \begin{cases} \emptyset, & i = 1, \dots, m, \\ g(X_{i-m}^i, Z_{i-m}^i), & i = m+1, m+2, \dots, \end{cases} \quad (21)$$

with the input X_i and the stationary ergodic noise process $\{Z_i\}_{i=1}^\infty$ depicted in Figure 1. We prove that the feedback capacity is given by

$$\begin{aligned} C_{FB} &= \lim_{n \rightarrow \infty} C_{FB,n} \\ &= \lim_{n \rightarrow \infty} \sup_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n), \end{aligned} \quad (22)$$

where the supremum is over all causally conditioned distributions

$$p(x^n || y^{n-1}) = \prod_{i=1}^n p(x_i | x^{i-1}, y^{i-1}).$$

We will combine the coding technique developed in the previous section with the Shannon strategy for channels with side information, in particular, Lemma 10.

That the limit in (22) is well-defined follows from the superadditivity of $nC_{FB,n}$. Thus,

$$C_{FB} = \lim_{n \rightarrow \infty} C_{FB,n} = \sup_{n \geq 1} C_{FB,n}.$$

The converse was proved by Massey [26, Theorem 3]. We repeat the proof here for completeness. For any sequence of $(2^{nR}, n)$ codes with $P_e^{(n)}$, we have from Fano's inequality

$$\begin{aligned} nR &\leq I(W; Y^n) + n\epsilon_n \\ &= \sum_{i=1}^n I(W; Y_i | Y^{i-1}) + n\epsilon_n \\ &= \sum_{i=1}^n I(X^i; Y_i | Y^{i-1}) + n\epsilon_n \\ &= I(X^n \rightarrow Y^n) + n\epsilon_n, \end{aligned} \quad (23)$$

where $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$. Here (23) follows from the codebook structure $X_i(W, Y^{i-1})$ and the Markovity $W \rightarrow (X^i, Y^{i-1}) \rightarrow Y_i$.

For the achievability, we show that there exists a sequence of codes that achieves $C_{FB,n}$ for each n . As before, we assume that the alphabets are finite. In the light of Lemma 10, it suffices to show that

$$C'_{FB,n} = \max_{p(u^n), x_i = f(u_i, x^{i-1}, y^{i-1})} I(U^n; Y^n) \quad (24)$$

is achievable, where the auxiliary random variables U_i has the cardinality bounded by $|\mathcal{U}_i| \leq |\mathcal{X}|^i |\mathcal{Y}|^{i-1}$, and the maximization is over all joint distributions of the form

$$p(u^n, x^n, y^n) = \left(\prod_{i=1}^n p(u_i) p(x_i | u_i, x^{i-1}, y^{i-1}) \right) p(y^n || x^n)$$

with deterministic $p(x_i | u_i, x^{i-1}, y^{i-1})$, $i = 1, \dots, n$.

Codebook generation and encoding. Fix n and let $p_i^*(u_i)$, $i = 1, \dots, n$, and $f_i^* : (u_i, x^{i-1}, y^{i-1}) \mapsto x_i$, $i = 1, \dots, n$, achieve the maximum of (24). We will also use the notation $p^*(u^n) = \prod_{i=1}^n p_i^*(u_i)$

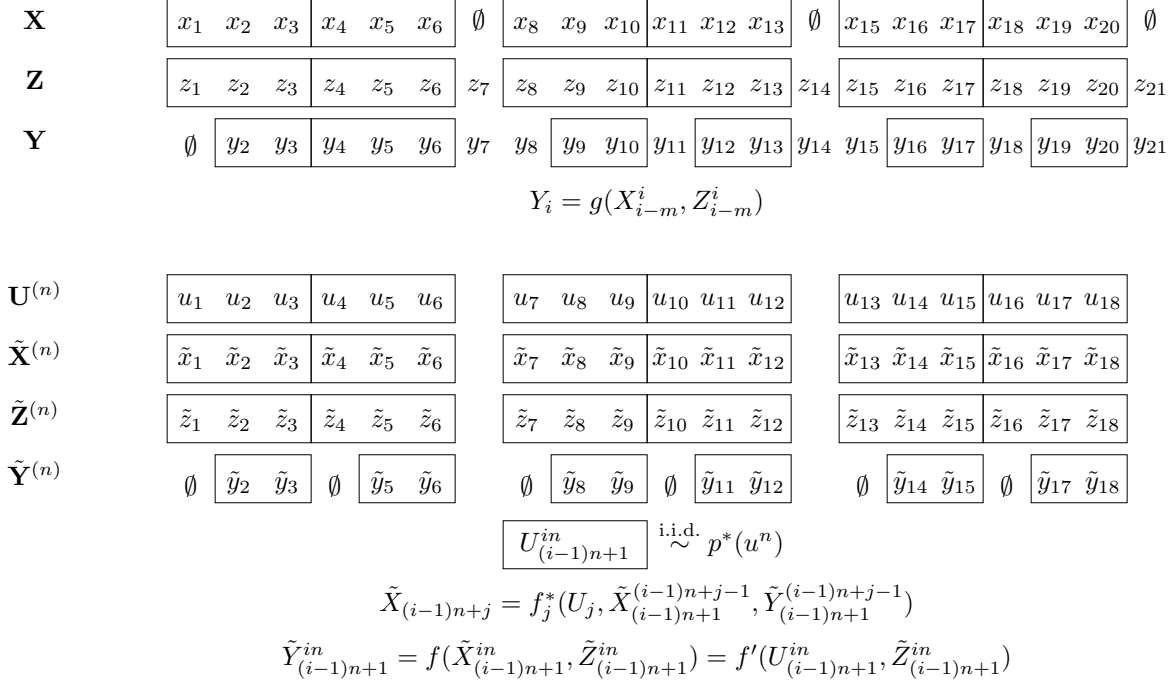


Figure 5: Code, input, noise, and output sequences: $n = 3, L = 2, m = 1$.

and $f^*(u^n, x^{n-1}, y^{n-1}) = (f_1^*(u_1), \dots, f_n^*(u_n, x^{n-1}, y^{n-1}))$.

For each $k = k(L, n) = Ln^2 + n$, $L = 1, 2, \dots$, we generate a $(2^{kR}, k)$ code $\{X_i(W, Y^{i-1})\}_{i=1}^k$ as summarized in Figure 5. As before, $\tilde{\mathbf{X}}^{(n)}$, $\tilde{\mathbf{Y}}^{(n)}$, and $\tilde{\mathbf{Z}}^{(n)}$ are respectively related to the underlying sequences $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ with every $(Ln + 1)$ st symbol omitted.

For each $w \in \{1, 2, \dots, 2^{kR}\}$, we generate a codeword $\mathbf{U}^{(n)}(w) = U^{Ln^2}(w)$ of length Ln on the n -letter alphabet $\mathcal{U}_1 \times \dots \times \mathcal{U}^n$ independently according to

$$p(u^{Ln^2}) = \prod_{i=1}^{Ln} p^*(u_{(n-1)i+1}^{ni}).$$

This gives a $2^{kR} \times Ln$ codebook matrix with each entry drawn i.i.d. according to $p^*(u^n)$.

To communicate the message $W = w$, the transmitter chooses the codeword $\mathbf{U}^{(n)}(w) = U^{Ln^2}(w)$ and sends

$$\tilde{X}_{(i-1)n+j} = f_j^*(U_j(w), \tilde{X}_{(i-1)n+1}^{(i-1)n+j-1}, \tilde{Y}_{(i-1)n+1}^{(i-1)n+j-1}), \quad i = 1, \dots, Ln, \quad j = 1, \dots, n.$$

Thus, the code function $X^n(w, Y^{n-1})$ utilizes the codeword $\mathbf{U}^{(n)}$ and the channel feedback $\tilde{\mathbf{Y}}^{(n)}$ only within the frame of n transmissions (each box in Figure 5).

Decoding. Upon receiving Y^k , the receiver declares that the message \hat{W} was sent if there is a unique \hat{W} such that

$$(\mathbf{U}^{(n)}(\hat{W}), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}(U^n, Y^n),$$

that is, $(\mathbf{U}^{(n)}(\hat{W}), \tilde{\mathbf{Y}}^{(n)})$ is jointly typical with respect to the joint distribution $p(u^n, y^n)$ specified

by $p^*(u^n)p(z^n)$, $x_i = f_i^*(u_i, x^{i-1}, y^{i-1})$, and the definition of the channel (21). Otherwise, an error is declared.

Analysis of the probability of error. We define the following events:

$$E_i = \{(\mathbf{U}^{(n)}(i), \tilde{\mathbf{Y}}^{(n)}) \in A_\epsilon^{*(Ln)}(U^n, Y^n)\}, \quad i \in \{1, 2, \dots, 2^{kR}\}.$$

As before, we assume $W = 1$ was sent.

From Lemma 7, $\mathbf{U}^{(n)}(1)$ and $\mathbf{Z}^{(n)}$ are jointly typical with high probability for L sufficiently large. Furthermore, $\tilde{\mathbf{Y}}^{(n)}$ is an n -letter blockwise function of $(\tilde{\mathbf{X}}^{(n)}(1), \tilde{\mathbf{Z}}^{(n)})$, and thus of $(\mathbf{U}^{(n)}(1), \tilde{\mathbf{Z}}^{(n)})$. Therefore, the probability of the event E_1^c that the intended codeword $\mathbf{U}^{(n)}(1)$ is not jointly typical with $\tilde{\mathbf{Y}}^{(n)}$ vanishes as $L \rightarrow \infty$.

On the other hand, $\mathbf{U}^{(n)}(i)$, $i \neq 1$, is generated blockwise i.i.d. $\sim p^*(u^n)$ independent of $\mathbf{Y}^{(n)}$. Hence, from Lemma 8, the probability of the event E_i that $\mathbf{U}^{(n)}(i)$ is jointly typical with $\mathbf{Y}^{(n)}$ is bounded by

$$\Pr(E_i) \leq 2^{-Ln(I(U^n; Y^n) - \delta)}, \quad \text{for all } i \neq 1,$$

where $\delta \rightarrow 0$ as $\epsilon \rightarrow 0$. Consequently, we have

$$\begin{aligned} \Pr(\hat{W} \neq W) &\leq \Pr(E_1^c) + \sum_{i=2}^{2^{kR}} \Pr(E_i) \\ &\leq \epsilon + 2^{kR} 2^{-Ln(I(U^n; Y^n) - \delta)} \\ &\leq 2\epsilon \end{aligned}$$

if L is sufficiently large and

$$kR = (Ln^2 + n)R < Ln(I(U^n; Y^n) - \delta).$$

Thus by letting $L \rightarrow \infty$ and then $\epsilon \rightarrow 0$, we can achieve any rate $R < C'_{FB,n}$.

Finally by Lemma 10, this implies that we can achieve

$$C_{FB,n} = \max_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n),$$

which completes the proof of Theorem 1.

5 Concluding Remarks

Trading off generality off for transparency, we have focused on the stationary channels of the form

$$Y_n = f(X_{n-m}^n, Z_{n-m}^n)$$

and presented a simple and constructive proof of the feedback coding theorem. The Shannon strategy (Lemma 10) has a fundamental role in transforming the feedback coding problem into a nonfeedback one, which is then solved by a scalable coding scheme of constructing a long typical input-output sequence pair by concatenating shorter nonergodic ones with appropriate phase shifts.

This two-stage approach can be applied to other channel models and give a straightforward coding theorem. For example, we can show that the finite-state channel

$$p(y_n, s_n | s_{n-1}, x_n) = p(y_n | s_{n-1}, x_n) p(s_n | s_{n-1}, x_n, y_n)$$

with deterministic $p(s_n | s_{n-1}, x_n, y_n)$ (but no assumption of indecomposability) has the feedback capacity lower bounded by

$$C_{FB} \geq \sup_{n \geq 1} \max_{p(x^n || y^{n-1})} \min_{s_0} \frac{1}{n} I(X^n \rightarrow Y^n | s_0).$$

This result was previously shown by Permuter *et al.* [31, Section V] via a generalization of Gallager's random coding exponent method for finite state channels without feedback [14, Section 5.9]. Here we sketch a simple alternative proof.

From a trivial modification of Lemma 10, the problem reduces to showing that

$$\max_{p(u^n), x_i = f(u_i, x^{i-1}, y^{i-1})} \min_{s_0} \frac{1}{n} I(U^n; Y^n | s_0) \quad (25)$$

is achievable for each n . But the given Shannon strategy ($p^*(u^n), x^n = f^*(u^n, x^{n-1}, y^{n-1})$) induces a new time-invariant finite-state channel on the n -letter super alphabet as $p(\mathbf{y}_k, \mathbf{s}_k | \mathbf{s}_{k-1}, \mathbf{u}_k)$. Hence we can use Gallager's random coding exponent method directly to achieve

$$\lim_{k \rightarrow \infty} \max_{p(\mathbf{u}^k)} \min_{s_0} \frac{1}{k} I(\mathbf{U}^k; \mathbf{Y}^k | s_0),$$

which can be shown to be larger than our target

$$\frac{1}{n} I(\mathbf{U}_1; \mathbf{Y}_1 | s_0),$$

because of the deterministic evolution of the state $S_n = f(S_{n-1}, X_n, Y_n)$.

We finally mention an important question that is not dealt with in this paper. Our characterization of the feedback capacity

$$C_{FB} = \lim_{n \rightarrow \infty} \max_{p(x^n || y^{n-1})} \frac{1}{n} I(X^n \rightarrow Y^n) \quad (26)$$

or any similar multi-letter expressions are in general not computable and do not provide much insight on the structure of the capacity achieving coding scheme. One may ask whether a stationary or even Markov distribution is asymptotically optimal for the sequence of maximizations in (26). This problem has been solved for a few specific channel models such as certain classes of finite-state channels [6, 40, 31, 30] and stationary additive Gaussian noise channels [21, 22], sometimes with analytic expressions for the feedback capacity. In this context, the current development is just the first step toward the complete characterization of the feedback capacity.

Acknowledgment

The author wishes to thank Tom Cover, Bob Gray, and Haim Permuter for helpful discussions.

References

- [1] F. Alajaji, “Feedback does not increase the capacity of discrete channels with additive noise,” *IEEE Trans. Inf. Theory*, vol. IT-41, no. 2, pp. 546–549, Mar. 1995.
- [2] T. Berger, *Rate Distortion Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1971.
- [3] D. Blackwell, “Information theory,” in *Modern Mathematics for the Engineer: Second Series*. New York: McGraw-Hill, 1961, pp. 182–193.
- [4] J. R. Brown, *Ergodic Theory and Topological Dynamics*. New York: Academic Press, 1976.
- [5] G. Caire and S. Shamai, “On the capacity of some channels with channel state information,” *IEEE Trans. Inf. Theory*, vol. IT-45, no. 6, pp. 2007–2019, 1999.
- [6] J. Chen and T. Berger, “The capacity of finite-state Markov channels with feedback,” *IEEE Trans. Inf. Theory*, vol. IT-51, no. 3, pp. 780–798, Mar. 2005.
- [7] T. M. Cover, “An achievable rate region for the broadcast channel,” *IEEE Trans. Inf. Theory*, vol. IT-21, pp. 399–404, 1975.
- [8] T. M. Cover and S. Pombra, “Gaussian feedback capacity,” *IEEE Trans. Inf. Theory*, vol. IT-35, no. 1, pp. 37–43, Jan. 1989.
- [9] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New York: Wiley, 2006.
- [10] R. L. Dobrushin, “General formulation of Shannon’s main theorem in information theory,” *Uspekhi Mat. Nauk*, vol. 14, no. 6, pp. 3–104, 1959, English transl. in *Amer. Math. Soc. Transl.*, vol. 33, no. 2, pp. 323–438, 1963.
- [11] A. Feinstein, “On the coding theorem and its converse for finite-memory channels,” *Information and Control*, vol. 2, pp. 25–44, 1959.
- [12] —, “A new basic theorem of information theory,” *IRE Trans. Inf. Theory*, vol. IT-4, pp. 2–22, 1954.
- [13] G. D. Forney, Jr., *Information Theory*, unpublished course notes, Stanford University, 1972.
- [14] R. G. Gallager, *Information Theory and Reliable Communication*. New York: Wiley, 1968.
- [15] —, “A simple derivation of the coding theorem and some applications,” *IEEE Trans. Inf. Theory*, vol. IT-11, pp. 3–18, 1965.
- [16] R. M. Gray, *Entropy and Information Theory*. New York: Springer-Verlag, 1990.

- [17] R. M. Gray and D. S. Ornstein, “Block coding for discrete stationary \bar{d} -continuous noisy channels,” *IEEE Trans. Inf. Theory*, vol. IT-25, no. 3, pp. 292–306, 1979.
- [18] T. S. Han, *Information-Spectrum Methods in Information Theory*. New York: Springer, 2003.
- [19] A. I. Khinchin, *Mathematical Foundations of Information Theory*. New York: Dover, 1957.
- [20] J. C. Kieffer, “Block coding for weakly continuous channels,” *IEEE Trans. Inf. Theory*, vol. IT-27, no. 6, pp. 721–727, 1981.
- [21] —, “Feedback capacity of the first-order moving average Gaussian channel,” *IEEE Trans. Inf. Theory*, vol. IT-52, no. 7, pp. 3063–3079, 2006.
- [22] Y.-H. Kim, “Feedback capacity of stationary Gaussian channels,” submitted to *IEEE Trans. Inf. Theory*, February 2006. [Online]. Available: <http://arxiv.org/abs/cs.IT/0602091/>
- [23] G. Kramer, *Directed Information for Channels with Feedback*. Konstanz: Hartung-Gorre Verlag, 1998, Dr. sc. thchn. Dissertation, Swiss Federal Institute of Technology (ETH) Zurich.
- [24] —, “Capacity results for the discrete memoryless network,” *IEEE Trans. Inf. Theory*, vol. IT-49, no. 1, pp. 4–21, 2003.
- [25] A. Lapidoth and İ. E. Telatar, “The compound channel capacity of a class of finite-state channels,” *IEEE Trans. Inf. Theory*, vol. IT-44, no. 3, pp. 973–983, 1998.
- [26] J. L. Massey, “Causality, feedback, and directed information,” in *Proc. International Symposium on Information Theory and its Applications*, Honolulu, Hawaii, Nov. 1990, pp. 303–305.
- [27] J. L. Massey and P. C. Massey, “Conservation of mutual and directed information,” in *Proc. International Symposium on Information Theory*, Adelaide, Australia, Sept. 2005, pp. 157–158.
- [28] J. Nedoma, “Über die Ergodizität und r -Ergodizität stationärer Wahrscheinlichkeitsmasse,” *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, vol. 2, pp. 90–97, 1963.
- [29] D. L. Neuhoff and P. C. Shields, “Channels with almost finite memory,” *IEEE Trans. Inf. Theory*, vol. IT-25, no. 4, pp. 440–447, 1979.
- [30] H. Permuter, P. Cuff, B. Van Roy, and T. Weissman, “Capacity of the trapdoor channel with feedback,” submitted to *IEEE Trans. Inform. Theory*, 2006. [Online]. Available: <http://arxiv.org/abs/cs.IT/0610047/>
- [31] H. Permuter, T. Weissman, and A. Goldsmith, “Finite-state channels with time-invariant deterministic feedback,” submitted to *IEEE Trans. Inform. Theory*, 2006. [Online]. Available: <http://arxiv.org/abs/cs.IT/0608070/>
- [32] K. Petersen, *Ergodic Theory*. Cambridge: Cambridge University Press, 1983.
- [33] M. S. Pinsker, *Information and Information Stability of Random Variables and Processes*. San Francisco: Holden-Day, 1964.

- [34] C. E. Shannon, “A mathematical theory of communication,” *Bell System Tech. J.*, vol. 27, pp. 379–423, 623–656, 1948.
- [35] —, “The zero error capacity of a noisy channel,” *IRE Trans. Inf. Theory*, vol. IT-2, no. 3, pp. 8–19, Sept. 1956.
- [36] —, “Channels with side information at the transmitter,” *IBM J. Res. Develop.*, vol. 2, pp. 289–293, 1958.
- [37] S. Tatikonda, “Control under communication constraints,” Ph.D. Thesis, Massachusetts Institute of Technology, Sept. 2000.
- [38] S. Verdú and T. S. Han, “A general formula for channel capacity,” *IEEE Trans. Inf. Theory*, vol. IT-40, no. 4, pp. 1147–1157, July 1994.
- [39] J. Wolfowitz, *Coding Theorems of Information Theory*, 3rd ed. Berlin: Springer-Verlag, 1978.
- [40] S. Yang, A. Kavčić, and S. Tatikonda, “Feedback capacity of finite-state machine channels,” *IEEE Trans. Inf. Theory*, vol. IT-51, no. 3, pp. 799–810, Mar. 2005.